## Chapter 5 Elements of Calculus

An Introduction to Optimization

Spring, 2014

Wei-Ta Chu

- A sequence of real numbers can be viewed as a set of numbers {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>k</sub>, ...}, which is often also denoted as {x<sub>k</sub>} or {x<sub>k</sub>}<sup>∞</sup><sub>k=1</sub>
- A sequence {x<sub>k</sub>} is *increasing* if x<sub>1</sub> < x<sub>2</sub> < ··· < x<sub>k</sub> < ···. If x<sub>k</sub> ≤ x<sub>k+1</sub> then we say that the sequence is *nondecreasing*. Similarly, we can define *decreasing* and *nonincreasing* sequences. Nonincreasing and nondecreasing sequences are called *monotone sequences*.
- A number x\* ∈ R is called the *limit* of the sequence {xk} if for any positive ε there is a number K (which may depend on ε) such that for all k > K, |xk x\*| < ε. In this case, we write x\* = lim<sub>k→∞</sub> xk or xk → x\*
- A sequence that has a limit is called a *convergent sequence*.

- A sequence in R<sup>n</sup> is a function whose domain is the set of natural numbers 1, 2, ..., k, ... and whose range is contained in R<sup>n</sup>. We use the notation {x<sup>(1)</sup>, x<sup>(2)</sup>, ...} or {x<sup>(k)</sup>} for sequences in R<sup>n</sup>.
- For limits of sequences in R<sup>n</sup>, we need to replace absolution values with vector norms. In other words, x\* is the limit of {x<sup>(k)</sup>} if for any positive ε there is a number K such that k > K, ||x<sup>(k)</sup> x\*|| < ε.</li>
- If a sequence  $\{x^{(k)}\}$  is convergent, we write  $x^* = \lim_{k \to \infty} x^{(k)}$  or  $x^{(k)} \to x^*$

- > Theorem 5.1: A convergent sequence has only one limit.
- A sequence {x<sup>(k)</sup>} in R<sup>n</sup> is bounded if there exists a number B ≥ 0 such that ||x<sup>(k)</sup>|| ≤ B for all k = 1, 2, ...
- Theorem 5.2: Every convergent sequence is bounded.
- For a sequence {x<sub>k</sub>} in R, a number B is called an upper bound if x<sub>k</sub> ≤ B for all k = 1, 2, .... In this case, we say {x<sub>k</sub>} is bounded above.
- A number B is called an *lower bound* if x<sub>k</sub> ≥ B for all k = 1, 2, ....
   In this case, we say {x<sub>k</sub>} is *bounded below*.

- Any sequence {x<sub>k</sub>} in R that has an upper bound has a *least upper bound* (also called the *supremum*), which is the smallest number B that is an upper bound of {x<sub>k</sub>}. Similarly, it has a *greatest lower bound* (also called *infimum*).
- If *B* is the least upper bound of the sequence {*x<sub>k</sub>*}, then *x<sub>k</sub>* ≤ *B* for all *k*, and for any *ϵ* > 0, there exists a number *K* such that *x<sub>K</sub>* > *B* − *ϵ*
- If *B* is the greatest lower bound of  $\{x_k\}$ , then  $x_k \ge B$  for all *k*, and for any  $\epsilon > 0$ , there exists a number *K* such that  $x_K < B + \epsilon$
- Theorem 5.3: Every monotone bounded sequence in *R* is convergent.

Given a sequence {x<sup>(k)</sup>} and an increasing sequence of natural numbers {m<sub>k</sub>}, the sequence

 $\{m{x}^{(m_k)}\} = \{m{x}^{(m_1)},m{x}^{(m_2)},...\}$ 

is called a *subsequence* of the sequence  $\{x^{(k)}\}$ .

- Theorem 5.4: Consider a convergent sequence {x<sup>(k)</sup>} with limit x\*. Then any subsequence of {x<sup>(k)</sup>} also converges to x\*.
- It turns out that any bounded sequence contains a convergent subsequence (*Bolzano-Weierstrass Theorem*)

Consider a function *f* : *R<sup>n</sup>* → *R<sup>m</sup>* and a point *x*<sub>0</sub> ∈ *R<sup>n</sup>*. Suppose that there exists *f*<sup>\*</sup> such that for any convergent sequence {*x*<sup>(k)</sup>} with limit *x*<sub>0</sub>, we have

$$\lim_{k o\infty}oldsymbol{f}(oldsymbol{x}^{(k)})=oldsymbol{f}^*$$

Then, we use the notation  $\lim_{x\to x_0} f(x)$  to represent  $f^*$ 

- It turns out that *f* is continuous at *x*<sub>0</sub> if and only if for any convergent sequence {*x*<sup>(k)</sup>} with limit *x*<sub>0</sub>, we have  $\lim_{k\to\infty} f(x^{(k)}) = f\left(\lim_{k\to\infty} x^{(k)}\right) = f(x_0)$
- Therefore, using the notation introduced above, the function *f* is continuous at *x*<sub>0</sub> if and only if lim<sub>*x*→*x*<sub>0</sub></sub> *f*(*x*) = *f*(*x*<sub>0</sub>)

- We say that a sequence {A<sub>k</sub>} of m × n matrices converges to the m × n matrix A if lim<sub>k→∞</sub> ||A − A<sub>k</sub>|| = 0.
- Lemma 5.1: Let A ∈ R<sup>n×n</sup>. Then, lim<sub>k→∞</sub> A<sup>k</sup> = O if and only if the eigenvalues of A satisfy |λ<sub>i</sub>(A)| < 1, i = 1, ..., n</p>
- Lemma 5.2: The series of  $n \times n$  matrices

 $\boldsymbol{I}_n + \boldsymbol{A} + \boldsymbol{A}^2 + \cdots + \boldsymbol{A}^k + \cdots$ 

converges if and only if  $\lim_{k\to\infty} A^k = O$ . In this case the sum of the series equals  $(I_n - A)^{-1}$ 

- A matrix-valued function A : R<sup>r</sup> → R<sup>n×n</sup> is continuous at a point ξ<sub>0</sub> ∈ R<sup>r</sup> if lim<sub>||ξ-ξ<sub>0</sub>||</sub> ||A(ξ) − A(ξ<sub>0</sub>)|| = 0
- Lemma 5.3: Let A : R<sup>r</sup> → R<sup>n×n</sup> be an n × n matrix-valued function that is continuous at ξ<sub>0</sub>. If A(ξ<sub>0</sub>)<sup>-1</sup>exists, then A(ξ)<sup>-1</sup> exists for ξ sufficiently close to ξ<sub>0</sub> and A(·)<sup>-1</sup> is continuous at ξ<sub>0</sub>.

# Differentiability

- Differential calculus is based on the idea of approximating an arbitrary function by an *affine function*.
- A function A : R<sup>n</sup> → R<sup>m</sup> is *affine* if there exists a *linear* function L : R<sup>n</sup> → R<sup>m</sup> and a vector y ∈ R<sup>m</sup> such that A(x) = L(x) + y

for every  $\boldsymbol{x} \in R^n$ 

10

- Consider a function *f* : R<sup>n</sup> → R<sup>m</sup> and a point *x*<sub>0</sub> ∈ R<sup>n</sup>. We wish to find an affine function A that approximates *f* near the point *x*<sub>0</sub>
- First, it is natural to impose the condition

$$\mathcal{A}(\boldsymbol{x}_0) = \boldsymbol{f}(\boldsymbol{x}_0)$$

 $\mathcal{A}(oldsymbol{x}_0) = oldsymbol{f}(oldsymbol{x}_0)$ 

# Differentiability

- ▶ Because  $\mathcal{A}(\boldsymbol{x}) = \mathcal{L}(\boldsymbol{x}) + \boldsymbol{y}$ , we obtain  $\boldsymbol{y} = \boldsymbol{f}(\boldsymbol{x}_0) \mathcal{L}(\boldsymbol{x}_0)$
- ▶ By the linearity of *L*,

$$\mathcal{L}(\boldsymbol{x}) + \boldsymbol{y} = \mathcal{L}(\boldsymbol{x}) - \mathcal{L}(\boldsymbol{x}_0) + \boldsymbol{f}(\boldsymbol{x}_0) = \mathcal{L}(\boldsymbol{x} - \boldsymbol{x}_0) + \boldsymbol{f}(\boldsymbol{x}_0)$$

• Hence, we may write

 $\mathcal{A}(\boldsymbol{x}) = \mathcal{L}(\boldsymbol{x} - \boldsymbol{x}_0) + \boldsymbol{f}(\boldsymbol{x}_0)$ 

Next, we require that A(x) approaches f(x) faster than x approaches x<sub>0</sub>; that is,

$$\lim_{\boldsymbol{x} \to \boldsymbol{x}_0, \boldsymbol{x} \in \Omega} \frac{\|\boldsymbol{f}(\boldsymbol{x}) - \mathcal{A}(\boldsymbol{x})\|}{\|\boldsymbol{x} - \boldsymbol{x}_0\|} = 0$$

The conditions ensure that A approximates f near x<sub>0</sub> in the sense that the approximation error is "small" compared with the distance of the point from x<sub>0</sub>.



$$\mathcal{A}(\boldsymbol{x}) = \mathcal{L}(\boldsymbol{x} - \boldsymbol{x}_0) + \boldsymbol{f}(\boldsymbol{x}_0) \quad \lim_{\boldsymbol{x} \to \boldsymbol{x}_0, \boldsymbol{x} \in \Omega} \frac{\|\boldsymbol{f}(\boldsymbol{x}) - \mathcal{A}(\boldsymbol{x})\|}{\|\boldsymbol{x} - \boldsymbol{x}_0\|} = 0$$
  
Differentiability

In summary, a function *f* : Ω → R<sup>m</sup>, Ω ⊂ R<sup>n</sup>, is said to be *differentiable* at *x*<sub>0</sub> ∈ Ω if there is an affine function that approximates *f* near *x*<sub>0</sub>; that is, there exists a linear function L : R<sup>n</sup> → R<sup>m</sup> such that

$$\lim_{\boldsymbol{x}\to\boldsymbol{x}_0,\boldsymbol{x}\in\Omega}\frac{\|\boldsymbol{f}(\boldsymbol{x})-(\boldsymbol{\mathcal{L}}(\boldsymbol{x}-\boldsymbol{x}_0)+\boldsymbol{f}(\boldsymbol{x}_0))\|}{\|\boldsymbol{x}-\boldsymbol{x}_0\|}=0$$

- The linear function L is determined uniquely by f and x<sub>0</sub> and is called the *derivative* of f at x<sub>0</sub>.
- The function is said to be *differentiable* on  $\Omega$  if f is differentiable at every point of its domain  $\Omega$ .



## Differentiability

In *R*, an affine function has the form *ax* + *b*, with *a, b* ∈ *R*.
 Hence, a real-valued function *f*(*x*) of a real variable *x* that is differentiable at *x*<sub>0</sub> can be approximated by a function

$$\mathcal{A}(x) = ax + b$$

• Because  $f(x_0) = \mathcal{A}(x_0) = ax_0 + b$ , we obtain

$$\mathcal{A}(x) = ax + b = a(x - x_0) + f(x_0)$$

The linear part of A(x), denoted by L(x) earlier, is just ax. The norm of a real number is its absolute value, so by the definition of differentiability

$$\lim_{x \to x_0} \frac{|f(x) - (a(x - x_0) + f(x_0))|}{|x - x_0|} = 0$$
  
$$\implies \lim_{x \to x_0} \frac{|f(x) - f(x_0)|}{|x - x_0|} = a$$

$$\lim_{x \to x_0} \frac{|f(x) - f(x_0)|}{|x - x_0|} = a \qquad \mathcal{A}(x) = ax + b = a(x - x_0) + f(x_0)$$
  
Differentiability

- The number *a* is commonly denoted  $f'(x_0)$  and is called the derivative of *f* at  $x_0$ .
- The affine function  $\mathcal{A}$  is therefore given by  $\mathcal{A}(x) = f(x_0) + f'(x_0)(x x_0)$
- The affine function is tangent to f at  $x_0$ .



$$\lim_{\boldsymbol{x}\to\boldsymbol{x}_0,\boldsymbol{x}\in\Omega}\frac{\|\boldsymbol{f}(\boldsymbol{x})-(\mathcal{L}(\boldsymbol{x}-\boldsymbol{x}_0)+\boldsymbol{f}(\boldsymbol{x}_0))\|}{\|\boldsymbol{x}-\boldsymbol{x}_0\|}=0$$

- Any linear transformation from R<sup>n</sup> to R<sup>m</sup>, and in particular the derivative L of f : R<sup>n</sup> → R<sup>m</sup>, can be represented by an m × n matrix.
- ▶ To find the matrix *L* of the derivative *L*, we use the natural basis {*e*<sub>1</sub>, *e*<sub>2</sub>, ..., *e<sub>n</sub>*} for *R<sup>n</sup>*. Consider the vectors

 $x_j = x_0 + te_j, j = 1, 2, ..., n$ 

By the definition of the derivative, we have

$$\lim_{t\to 0} \frac{f(x_j) - (tLe_j + f(x_0))}{t} = 0 \qquad j = 1, 2, .., n$$

This means that

$$\lim_{t\to 0} \frac{\boldsymbol{f}(\boldsymbol{x}_j) - \boldsymbol{f}(\boldsymbol{x}_0)}{t} = \boldsymbol{L}\boldsymbol{e}_j$$



The Derivative Matrix

$$\lim_{t\to 0} \frac{\boldsymbol{f}(\boldsymbol{x}_j) - \boldsymbol{f}(\boldsymbol{x}_0)}{t} = \boldsymbol{L}\boldsymbol{e}_j$$

- But  $Le_j$  is the *j*th column of the matrix L. The vector  $x_j$  differs from  $x_0$  only in the *j*th coordinate, and in that coordinate the difference is just the number *t*. Therefore, the left side is the partial derivative  $\frac{\partial f}{\partial x_i}(x_0)$
- Because vector limits are computed by taking the limit of each coordinate function, it follows that if

$$\boldsymbol{f}(\boldsymbol{x}) = \begin{bmatrix} f_1(\boldsymbol{x}) \\ \vdots \\ f_m(\boldsymbol{x}) \end{bmatrix} \quad \left( \boldsymbol{f}_i : R^n \to R, i = 1, ..., m \right)$$
  
then  $\frac{\partial \boldsymbol{f}}{\partial x_j}(\boldsymbol{x}_0) = \begin{bmatrix} \frac{\partial f_1}{\partial x_j}(\boldsymbol{x}_0) \\ \vdots \\ \frac{\partial f_m}{\partial x_j}(\boldsymbol{x}_0) \end{bmatrix}$  and the matrix  $\boldsymbol{L}$  has the form  
 $\begin{bmatrix} \frac{\partial f}{\partial x_1}(\boldsymbol{x}_0) \cdots \frac{\partial f}{\partial x_n}(\boldsymbol{x}_0) \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\boldsymbol{x}_0) & \cdots & \frac{\partial f_1}{\partial x_n}(\boldsymbol{x}_0) \\ \vdots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\boldsymbol{x}_0) & \cdots & \frac{\partial f_m}{\partial x_n}(\boldsymbol{x}_0) \end{bmatrix}$ 

### The Derivative Matrix

17

- The matrix L is called the *Jacobian matrix*, or *derivative matrix*, of f at  $x_0$ , and is denoted  $Df(x_0)$ .
- For convenience, we often refer to Df(x<sub>0</sub>) simply as the derivative of f at x<sub>0</sub>.



 $\mathcal{A}(x) = f(x_0) + f'(x_0)(x - x_0)$ 

## The Derivative Matrix

Theorem 5.5: If a function *f* : *R<sup>n</sup>* → *R<sup>m</sup>* is differentiable at *x*<sub>0</sub>, then the derivative of *f* at *x*<sub>0</sub> is determined uniquely and is represented by an *m* × *n* derivative matrix *Df*(*x*<sub>0</sub>). The best affine approximation of *f* near *x*<sub>0</sub> is then given by

$$\mathcal{A}(\boldsymbol{x}) = \boldsymbol{f}(\boldsymbol{x}_0) + D\boldsymbol{f}(\boldsymbol{x}_0)(\boldsymbol{x} - \boldsymbol{x}_0)$$

in the sense that

$$oldsymbol{f}(oldsymbol{x}) = \mathcal{A}(oldsymbol{x}) + oldsymbol{r}(oldsymbol{x})$$

and  $\lim_{x\to x_0} ||r(x)|| / ||x - x_0|| = 0$ .

The columns of the derivative matrix  $Df(x_0)$  are vector partial derivatives. The vector  $\frac{\partial f}{\partial x_j}(x_0)$  is a tangent vector at  $x_0$  to the curve f obtained by varying only the *j*th coordinate of x.

• If  $f: \mathbb{R}^n \to \mathbb{R}$  is differentiable, then the function  $\nabla f$  defined by  $\nabla f(\boldsymbol{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\boldsymbol{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\boldsymbol{x}) \end{bmatrix} = Df(\boldsymbol{x})^T$ 

is called the *gradient* of f.

Given f: R<sup>n</sup> → R, if ∇f is differentiable, we say that f is *twice* differentiable, and we write the derivative of ∇f as

$$D^{2}f = \begin{bmatrix} \frac{\partial^{2}f}{\partial x_{1}^{2}} & \frac{\partial^{2}f}{\partial x_{2}\partial x_{1}} & \cdots & \frac{\partial^{2}f}{\partial x_{n}\partial x_{1}} \\ \frac{\partial^{2}f}{\partial x_{1}\partial x_{2}} & \frac{\partial^{2}f}{\partial x_{2}^{2}} & \cdots & \frac{\partial^{2}f}{\partial x_{n}\partial x_{2}} \\ \vdots & & \ddots & \vdots \\ \frac{\partial^{2}f}{\partial x_{1}\partial x_{n}} & \frac{\partial^{2}f}{\partial x_{2}\partial x_{n}} & \cdots & \frac{\partial^{2}f}{\partial x_{n}^{2}} \end{bmatrix}$$

 $\frac{\partial^2 f}{\partial x_i \partial x_j}$  represents taking the partial derivative with respect to  $x_j$ first, then with respect to  $x_i$ 

## The Derivative Matrix

- The matrix D<sup>2</sup>f(x) is called the Hessian matrix of f at x, and is often also denoted F(x).
- A function *f* : Ω → R<sup>m</sup>, Ω ⊂ R<sup>n</sup>, is said to be *continuously differentiable* on Ω if it is differentiable (on Ω), and D*f* : Ω → R<sup>m×n</sup> is continuous; that is, the components of *f* have continuous partial derivatives. In this case, we write *f* ∈ C<sup>1</sup>. If the components of *f* have continuous partial derivatives of order *p*, we write *f* ∈ C<sup>p</sup>.
- The Hessian matrix of a function f : R<sup>n</sup> → R at x is symmetric if f is twice continuously differentiable at x. This is a well-known result from calculus called *Clairaut's theorem* or *Schwarz's theorem*.

Example  

$$\begin{bmatrix} \frac{d}{dx} \left( \frac{f}{g} \right) = \frac{f'g - fg'}{g^2} \end{bmatrix} \quad F = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$
• Consider the function  $f(\mathbf{x}) = \begin{cases} x_1 x_2 (x_1^2 - x_2^2) / (x_1^2 + x_2^2) & \text{if } \mathbf{x} \neq \mathbf{0} \\ 0 & \text{if } \mathbf{x} = \mathbf{0} \end{cases}$ 

Compute its Hessian at the point  $\mathbf{0} = [0, 0]^T$ 

Start with 
$$\frac{\partial^2 f}{\partial x_1^2} = \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_1} \right)$$
$$\frac{\partial f}{\partial x_1} (\boldsymbol{x}) = \begin{cases} x_2 (x_1^4 - x_2^4 + 4x_1^2 x_2^2) / (x_1^2 + x_2^2)^2 & \text{if } \boldsymbol{x} \neq \boldsymbol{0} \\ 0 & \text{if } \boldsymbol{x} = \boldsymbol{0} \end{cases}$$

Note that

$$\frac{\partial f}{\partial x_1}([x_1, 0]^T) = 0 \qquad \qquad \frac{\partial^2 f}{\partial x_1^2}(\mathbf{0}) = 0$$
$$\frac{\partial f}{\partial x_1}([0, x_2]^T) = -x_2 \qquad \qquad \frac{\partial f}{\partial x_2 x_1}(\mathbf{0}) = -1$$

### Example

• We next compute

$$\frac{\partial^2 f}{\partial x_2^2} = \frac{\partial}{\partial x_2} \left( \frac{\partial f}{\partial x_2} \right)$$

$$\frac{\partial f}{\partial x_2} (\boldsymbol{x}) = \begin{cases} x_1 (x_1^4 - x_2^4 - 4x_1^2 x_2^2) / (x_1^2 + x_2^2)^2 & \text{if } \boldsymbol{x} \neq \boldsymbol{0} \\ 0 & \text{if } \boldsymbol{x} = \boldsymbol{0} \end{cases}$$

$$\frac{\partial f}{\partial x_2} ([0, x_2]^T) = 0 \qquad \frac{\partial f}{\partial x_2^2} (\boldsymbol{0}) = 0$$

$$\frac{\partial f}{\partial x_2} ([x_1, 0]^T) = x_1 \qquad \frac{\partial f}{\partial x_1 x_2} ([x_1, 0]^T) = 1$$

• Therefore, the Hessian evaluated at the point  $\mathbf{0} = [0, 0]^T$  is  $\mathbf{F}(\mathbf{0}) = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ 

# **Differentiation Rules**



- The *chain rule* for differentiating the composition g(f(t)), of a function f: R → R<sup>n</sup> and a function g: R<sup>n</sup> → R
- Theorem 5.6: Let g: D → R be differentiable on an open set D ⊂ R<sup>n</sup>, and let f: (a, b) → D be differentiable on (a, b). Then, the composite function h: (a, b) → R given by h(t) = g(f(t)) is differentiable on (a, b), and

$$h'(t) = Dg(\boldsymbol{f}(t))D\boldsymbol{f}(t) = \nabla g(\boldsymbol{f}(t))^T \begin{bmatrix} f_1'(t) \\ \vdots \\ f_n'(t) \end{bmatrix}$$

# **Differentiation Rules**



- *Product rule*: Let  $f : R^n \to R^m$  and  $g : R^n \to R^m$  be two differentiable functions. Define the function  $h : R^n \to R$  by  $h(x) = f(x)^T g(x)$ . Then h is also differentiable and  $Dh(x) = f(x)^T Dg(x) + g(x)^T Df(x)$
- Some useful formulas: Let A ∈ R<sup>m×n</sup> and y ∈ R<sup>m</sup>, the derivative with respect to x

$$D(\boldsymbol{y}^{T}\boldsymbol{A}\boldsymbol{x}) = \boldsymbol{y}^{T}\boldsymbol{A}$$
$$D(\boldsymbol{x}^{T}\boldsymbol{A}\boldsymbol{x}) = \boldsymbol{x}^{T}(\boldsymbol{A} + \boldsymbol{A}^{T}), \quad \text{if } m = n$$

- If  $\boldsymbol{y} \in R^n$ , then  $D(\boldsymbol{y}^T \boldsymbol{x}) = \boldsymbol{y}^T$
- If Q is a symmetric matrix, then  $D(\mathbf{x}^T Q \mathbf{x}) = 2\mathbf{x}^T Q$ . In particular,  $D(\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}^T$

- The *level set* of a function f : R<sup>n</sup> → R at level c is the set of points S = {x : f(x) = c}. For f : R<sup>2</sup> → R, we are usually interested in S when it is a curve. For f : R<sup>3</sup> → R, the sets S most often considered are surfaces.
- Example: Consider  $f(\boldsymbol{x}) = 100(x_2 x_1^2)^2 + (1 x_1)^2, \boldsymbol{x} = [x_1, x_2]^T$ It is called *Rosenbrock's function*.



- A point x<sub>0</sub> is on the level set S at level c means that f(x<sub>0</sub>) = c. Suppose that there is a curve γ lying in S and parameterized by a continuously differentiable function g : R → R<sup>n</sup>. Suppose also that g(t<sub>0</sub>) = x<sub>0</sub> and Dg(t<sub>0</sub>) = v ≠ 0, so that v is a tangent vector to γ at x<sub>0</sub>.
- Applying the chain rule to the function h(t) = f(g(t))  $h'(t_0) = Df(g(t_0))Dg(t_0) = Df(x_0)v$ Since  $\gamma$  lies on S, we have h(t) = f(g(t)) = cThat is, h is constant. Thus,  $h'(t_0) = 0 = Df(x_0)v = \nabla f(x_0)^T v$



X1

- Theorem 5.7: The vector 
  ¬
  f(x<sub>0</sub>) is orthogonal to the tangent vector to an arbitrary smooth curve passing through x<sub>0</sub> on the level set determined by f(x) = f(x<sub>0</sub>)
- It is natural to say that 
   ¬f(x<sub>0</sub>) is *orthogonal* or *normal* to the level set S corresponding to x<sub>0</sub>, and to take as the tangent plane (or line) to S at x<sub>0</sub> the set of all points x satisfying
   ¬f(x<sub>0</sub>)<sup>T</sup>(x x<sub>0</sub>) = 0, if ¬f(x<sub>0</sub>) ≠ 0



- $\nabla f(\boldsymbol{x}_0)$  is the direction of *maximum rate of increase* of f at  $\boldsymbol{x}_0$
- The direction of maximum rate of increase of a real-valued differentiable function at a point is orthogonal to the level set of the function through that point.
- An example about  $f : R^2 \to R$ 
  - The curve on the surface running from bottom to top has the property that its projection onto the (x<sub>1</sub>, x<sub>2</sub>)-plane is always orthogonal to the level curves and is called a *path of steepest ascent*.





- The *graph* of  $f: \mathbb{R}^n \to \mathbb{R}$  is the set  $\{[\mathbf{x}^T, f(\mathbf{x})]^T : \mathbf{x} \in \mathbb{R}^n\} \subset \mathbb{R}^{n+1}$ . The notion of the gradient of a function has an alternative useful interpretation in terms of the tangent hyperplane to its graph.
- Let  $x_0 \in \mathbb{R}^n$  and  $z_0 = f(x_0)$ . The point  $[x_0^T, z_0]^T \in \mathbb{R}^{n+1}$  is a point on the graph of f. If f is differentiable at  $\boldsymbol{\xi}$ , then the graph admits a nonvertical tangent hyperplane at  $\boldsymbol{\xi} = [\boldsymbol{x}^T, z_0]^T$ . The hyperplane through  $\boldsymbol{\xi}$  is the set of all points  $[x_1, ..., x_n, z]^T \in \mathbb{R}^{n+1}$ satisfying the equation

 $u_1(x_1 - x_{01}) + \dots + u_n(x_n - x_{0n}) + v(z - z_0) = 0$ where the vector  $[u_1, \dots, u_n, v] \in \mathbb{R}^{n+1}$  is normal to the hyperplane.

$$\langle (u_1, ..., u_n, v), (x_1 - x_{01}, ..., x_n - x_{0n}, z - z_0) \rangle$$
  $x_0 = (x_{01}, ..., x_{0n})$   
normal a vector on the hyperplane

• Assuming that this hyperplane is nonvertical (that is,  $v \neq 0$ ), let

$$d_i = -\frac{u_i}{v}$$

Thus, we can rewrite the hyperplane equation as

$$z = d_1(x_1 - x_{01}) + \dots + d_n(x_n - x_{0n}) + z_0$$

We can think of the right side as a function z : R<sup>n</sup> → R. Observe that for the hyperplane to be tangent to the graph of f, the functions f and z must have the same partial derivatives at the point x<sub>0</sub>. Hence, if f is differential at x<sub>0</sub>, its tangent hyperplane can be written in terms of its gradient, as given by the equation

$$z - z_0 = Df(\boldsymbol{x}_0)(\boldsymbol{x} - \boldsymbol{x}_0) = (\boldsymbol{x} - \boldsymbol{x}_0)^T \bigtriangledown f(\boldsymbol{x}_0)$$
 $Df(\boldsymbol{x}_0) = rac{z - z_0}{\boldsymbol{x} - \boldsymbol{x}_0}$ 

Theorem 5.8 *Taylor's Theorem*: Assume that a function f : R → R is m times continuously differentiable (i.e. f ∈ C<sup>m</sup>) on an interval [a, b]. Denote h = b - a Then,

\_\_\_\_\_

$$f(b) = f(a) + \frac{h}{1!}f^{(1)}(a) + \frac{h^2}{2!}f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!}f^{(m-1)}(a) + R_m$$

(called Taylor's formula) where  $f^{(i)}$  is the *i*th derivative of f, and

$$R_m = \frac{h^m (1-\theta)^{m-1}}{(m-1)!} f^{(m)}(a+\theta h) = \frac{h^m}{m!} f^{(m)}(a+\theta' h) \qquad \theta, \theta' \in (0,1)$$

- An important property of Taylor's theorem arises from the forms of the remainder  $R_m$ .
- We introduce the *order symbols*, *O* and *o*.
- Let g be a real-valued function defined in some neighborhood of 0 ∈ R<sup>n</sup>, with g(x) ≠ 0 if x ≠ 0. Let f : Ω → R<sup>m</sup> be defined in a domain Ω ⊂ R<sup>n</sup> that includes 0. Then, we write
  - 1. f(x) = O(g(x)) to mean that the quotient ||f(x)||/|g(x) is bounded near 0; that is, these exist numbers K > 0 and δ > 0 such that if ||x|| < δ, x ∈ Ω, then ||f(x)||/|g(x)| ≤ K or ||f(x)|| ≤ Kg(x)</li>
  - ▶ 2. f(x) = o(g(x)) to mean that

$$\lim_{\boldsymbol{x}\to\boldsymbol{0},\boldsymbol{x}\in\Omega}\frac{\|\boldsymbol{f}(\boldsymbol{x})\|}{|g(\boldsymbol{x})|}=0$$

- The symbol O(g(x)) [read "big-oh" of g(x)] is used to represent a function that is bounded by a scaled version of g in a neighborhood of 0.
- Examples:

$$x = O(x)$$

$$\begin{bmatrix} x^3 \\ 2x^2 + 3x^4 \end{bmatrix} = O(x^2)$$

$$\blacktriangleright \quad \sin x = O(x)$$

On the other hand, o(g(x)) [read "little-oh" of g(x)] represents a function that goes to zero "faster" than g(x) in the sense that lim<sub>x→0,x∈Ω</sub> ||o(g(x))||/|g(x)| = 0

• Examples:

$$x^{2} = o(x)$$

$$\begin{bmatrix} x^{3} \\ 2x^{2} + 3x^{4} \end{bmatrix} = o(x)$$

$$x^{3} = o(x^{2})$$

$$x = o(1)$$

- Note that if f(x) = o(g(x)), then f(x) = O(g(x)) (but the converse is not necessarily true). Also, if f(x) = O(||x||<sup>p</sup>), then f(x) = o(||x||<sup>p-ϵ</sup>) for any ϵ > 0
- Suppose that f ∈ C<sup>m</sup>. Recall that the remainder term in Taylor's theorem has the form

$$R_m = \frac{h^m}{m!} f^{(m)}(a + \theta h) \qquad \theta \in (0, 1)$$

Substituting this into Taylor's formula, we get

$$f(b) = f(a) + \frac{h}{1!}f^{(1)}(a) + \frac{h^2}{2!}f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!}f^{(m-1)}(a) + \frac{h^m}{m!}f^{(m)}(a+\theta h)$$

By the continuity of  $f^{(m)}$ , we have  $f^{(m)}(a + \theta h) \rightarrow f^{(m)}(a)$  as  $h \rightarrow 0$  that is,  $f^{(m)}(a + \theta h) = f^{(m)}(a) + o(1)$ . Therefore,

$$R_m = \frac{h^m}{m!} f^{(m)}(a + \theta h) = \frac{h^m}{m!} f^{(m)}(a) + o(h^m) \qquad h^m o(1) = o(h^m)$$

# • We may then write Taylor's formula as $f(b) = f(a) + \frac{h}{1!}f^{(1)}(a) + \frac{h^2}{2!}f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!}f^{(m-1)}(a) + \frac{h^m}{m!}f^{(m)}(a) + o(h^m)$

If, in addition, we assume that f ∈ C<sup>m+1</sup>, we may replace the term o(h<sup>m</sup>) above by O(h<sup>m+1</sup>)

$$f(b) = f(a) + \frac{h}{1!}f^{(1)}(a) + \frac{h^2}{2!}f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!}f^{(m-1)}(a) + \frac{h^m}{m!}f^{(m)}(a) + O(h^{m+1})$$

- Theorem 5.9 Mean value theorem: If a function f : R<sup>n</sup> → R<sup>m</sup> is differentiable on an open set Ω ⊂ R<sup>n</sup>, then for any pair of points x, y ∈ Ω, there exists a matrix M such that f(x) f(y) = M(x y)
- The mean value theorem follows from Taylor's theorem (for the case where m = 1) applied to each component of f. M is a matrix whose rows are the rows of Df evaluated at points that lie one the line segment joining x and y.

